

**NRES\_798\_7\_201501**

Statistical methods and linear models

# Basic statistical model

$$Y = \text{deterministic part} + \text{stochastic part}$$



Univariate  
Multivariate



Linear  
Nonlinear  
Smoothed



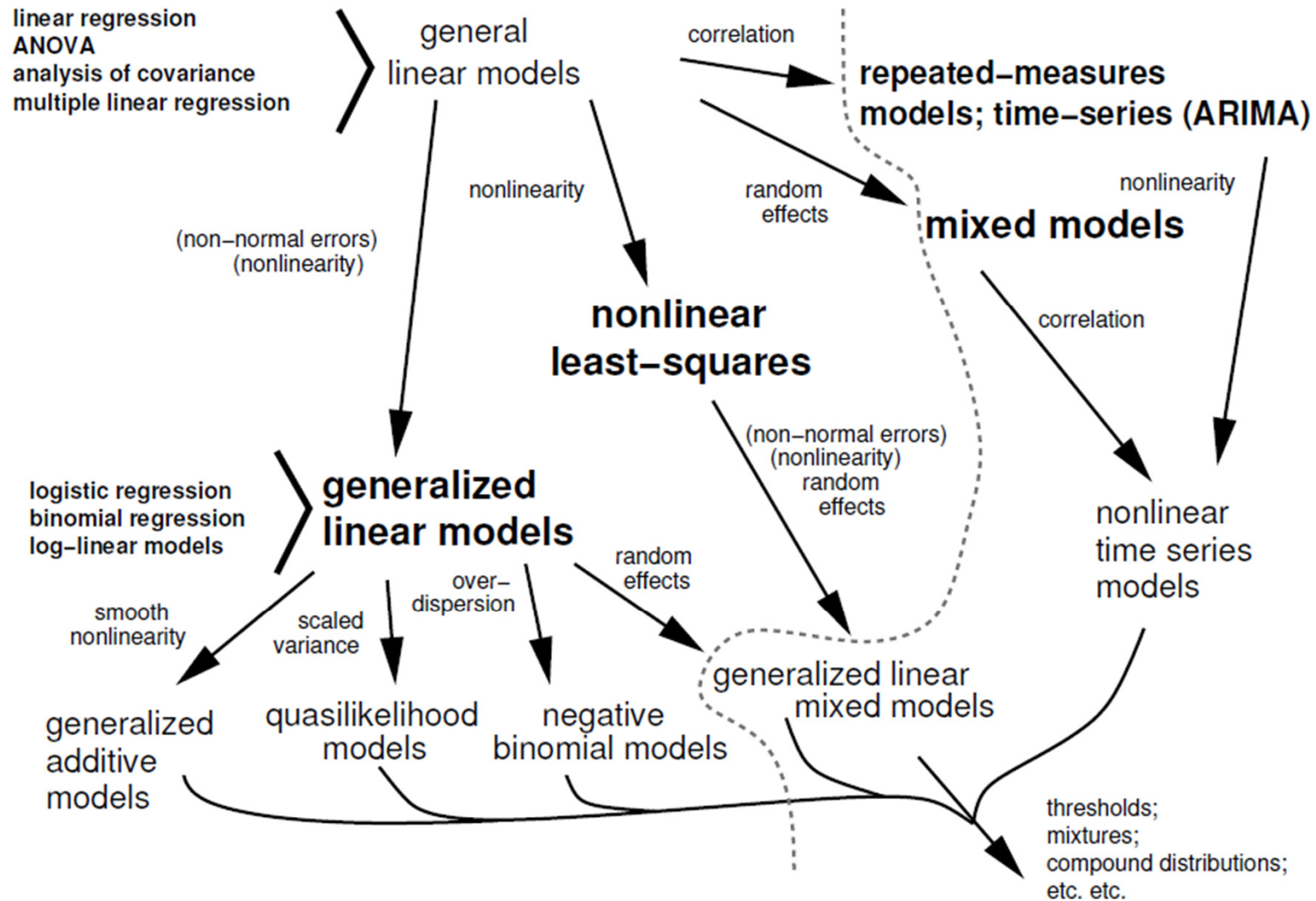
Distribution  
Heterogeneity  
Auto-correlation  
Nested data  
(random effects)  
Random noise

Error term ( $\varepsilon_i$ )

# General linear models

- Linear regression
- One- and multi-way analysis of variance (ANOVA)
- Analysis of covariance (ANCOVA)
- In R these procedures are use the function `lm()`

# Beyond linear models



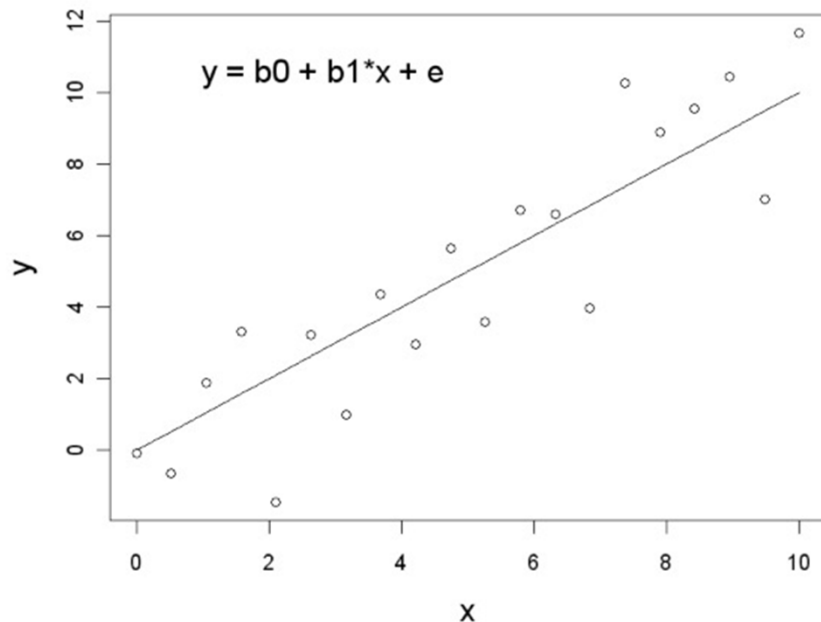
# General linear models

Deterministic  
part

$$Y \sim \text{Normal}(b_0 + b_1 X, \sigma^2)$$

Stochastic  
part

General linear model



Models that are linear functions of the parameters, not necessarily of the independent variables

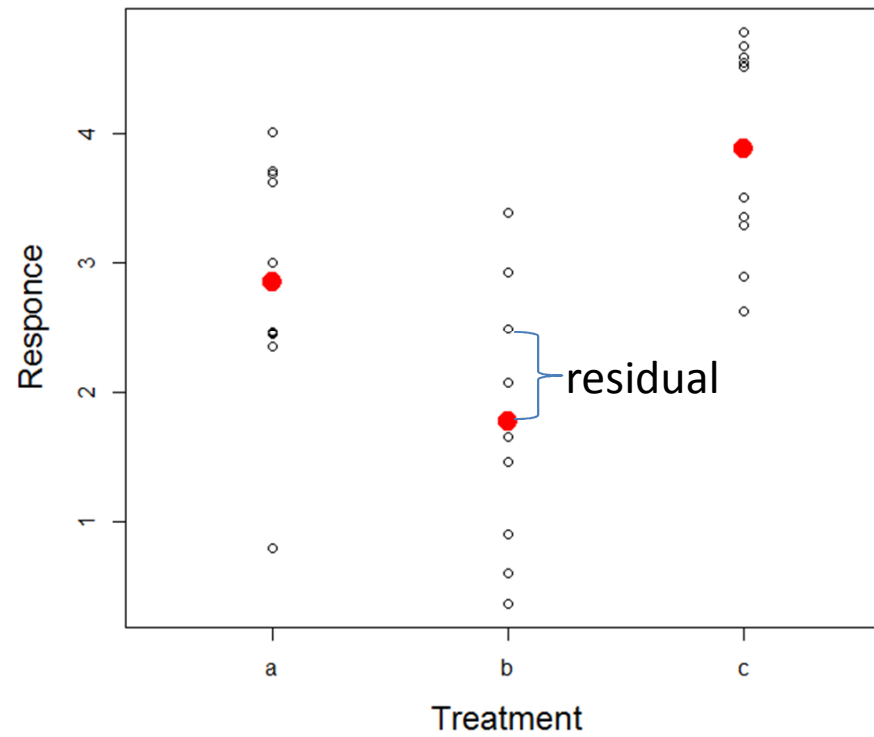
# General linear models

- Assumptions  $Y \sim \text{Normal}(b_0 + b_1X, \sigma^2)$ 
  - All observed values are:
    - Independent
    - Any continuous predictor variables (covariates) are measured without error
    - Constant variance (homoscedastic)
    - Normally distributed

# Homoscedasticity

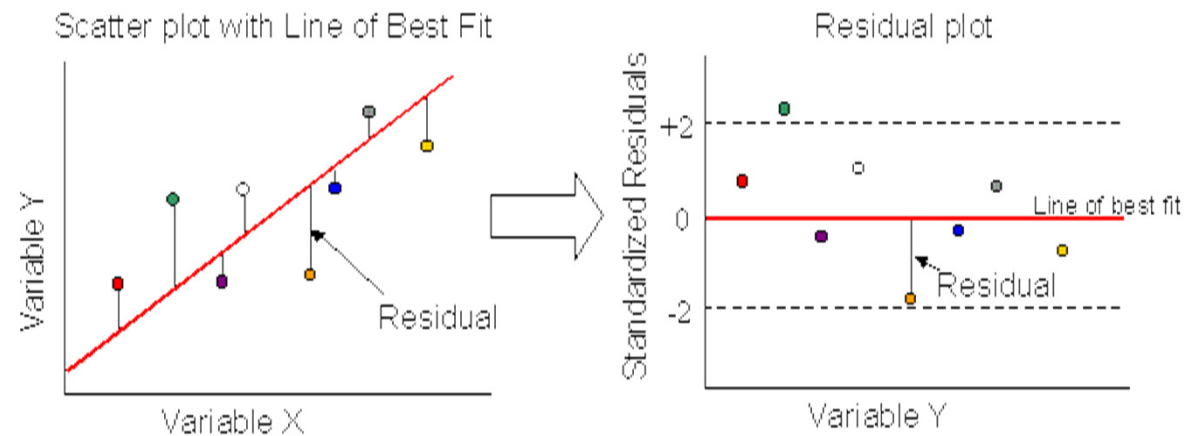
- Variance of all treatment groups needs to be approximately equal to each other
- Variance needs to be constant across all predictor variables
- Residuals

# ANOVA residuals

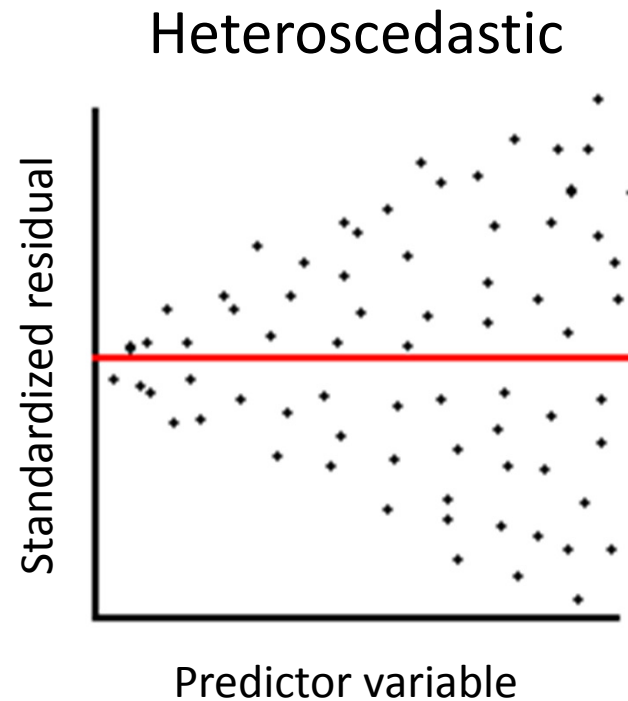
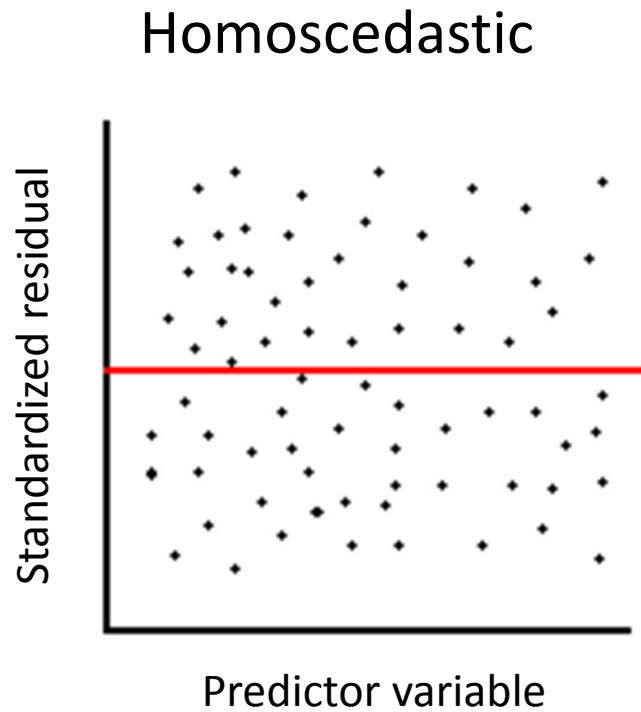




# Regression residuals



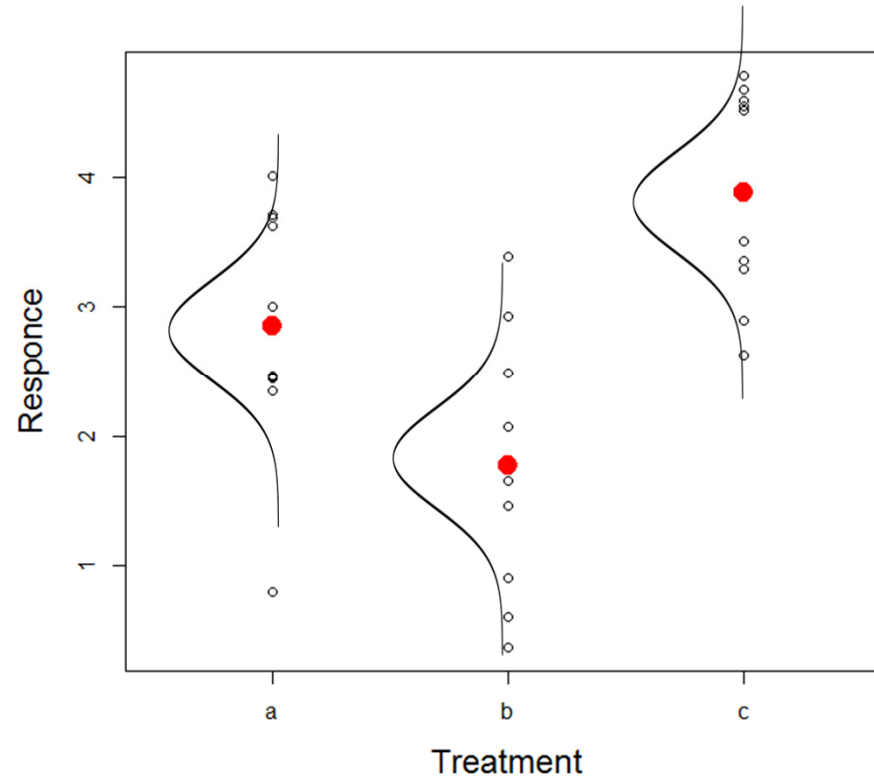
# Homoscedasticity



# Normality assumption

The assumption of normality applies to the variation around the expected value – the residual – not to the whole data set

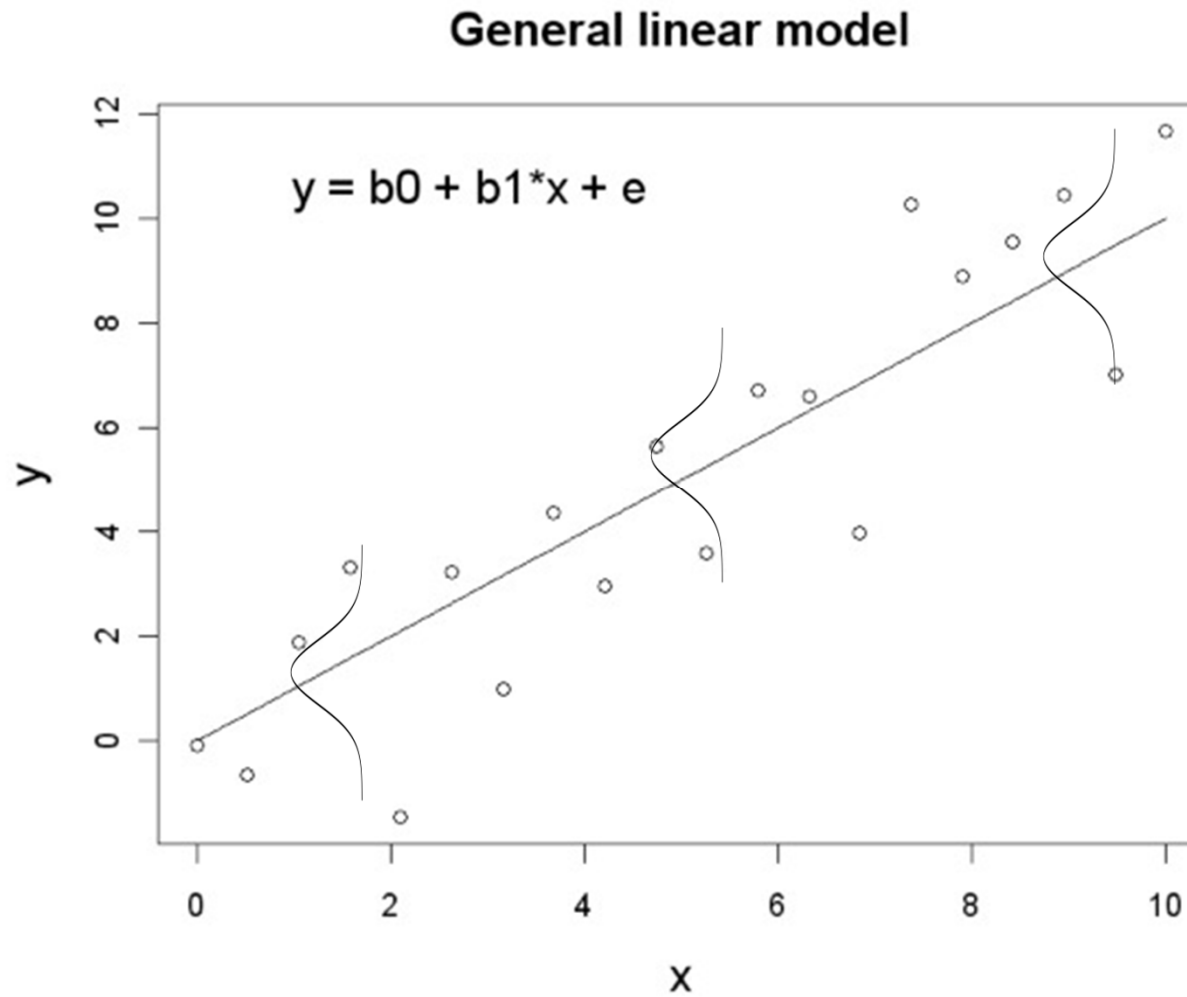
# Normality

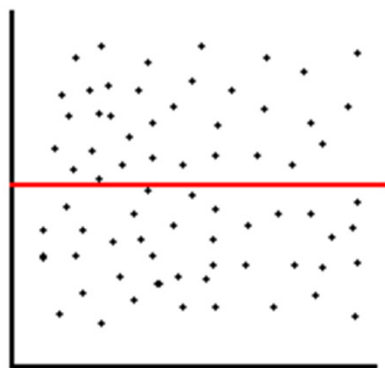


testing for ANOVA normality is easy

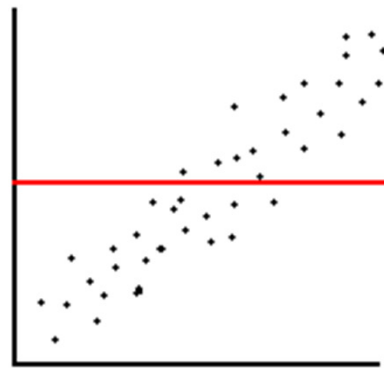
# Normality

$$Y \sim \text{Normal}(b_0 + b_1X, \sigma^2)$$

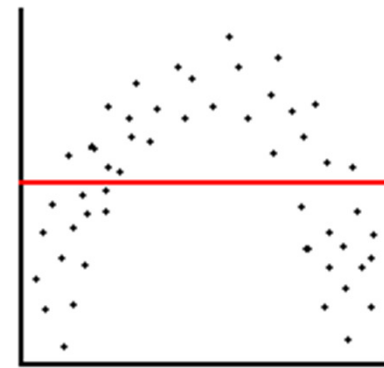




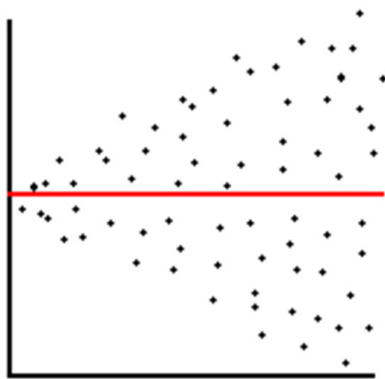
(a) Unbiased and Homoscedastic



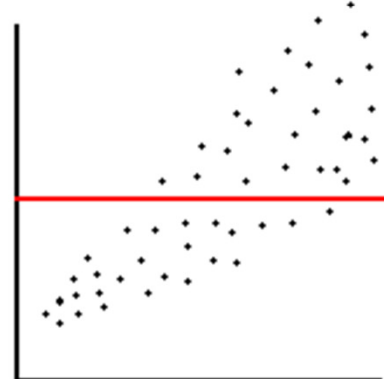
(b) Biased and Homoscedastic



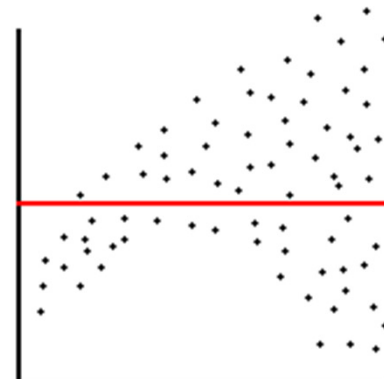
(c) Biased and Homoscedastic



(d) Unbiased and Heteroscedastic

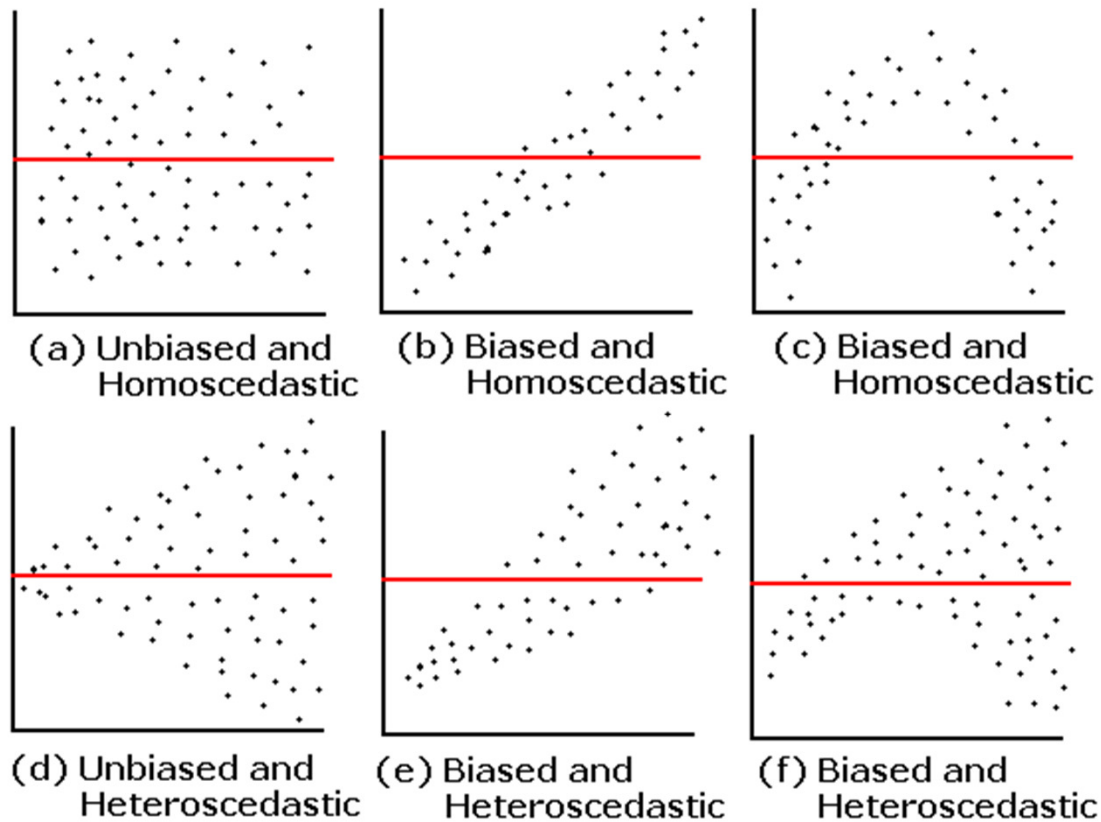


(e) Biased and Heteroscedastic



(f) Biased and Heteroscedastic

What do these residuals tell us about the statistical model?



- (a) Unbiased and homoscedastic. The residuals average to zero in each thin vertical strip and the SD is the same all across the plot.
- (b) Biased and homoscedastic. The residuals show a linear pattern, probably due to a lurking variable not included in the experiment.
- (c) Biased and homoscedastic. The residuals show a quadratic pattern, possibly because of a nonlinear relationship. Sometimes a variable transform will eliminate the bias.
- (d) Unbiased, but homoscedastic. The SD is small to the left of the plot and large to the right: the residuals are heteroscedastic.
- (e) Biased and heteroscedastic. The pattern is linear.
- (f) Biased and heteroscedastic. The pattern is quadratic.

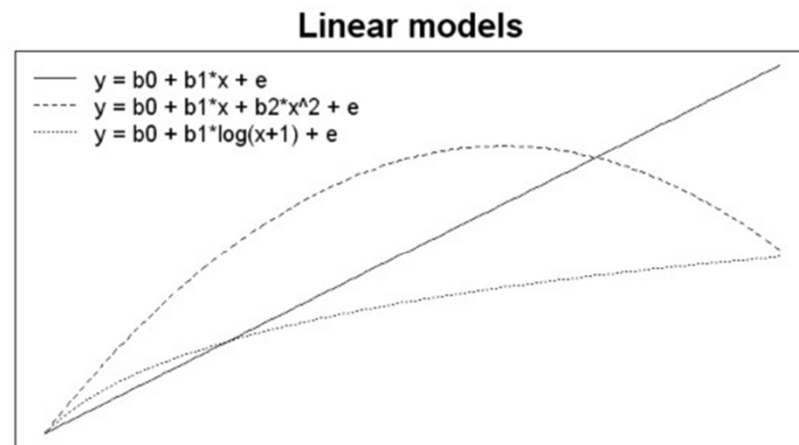
# Linear models

Example linear models:

$$Y \sim \text{Normal}(b_0 + b_1x, \sigma^2)$$

$$Y \sim \text{Normal}(b_0 + b_1x + b_2x^2, \sigma^2)$$

$$Y \sim \text{Normal}(b_0 + b_1 \log(x), \sigma^2)$$



In all of these models we can define a new explanatory variable  $Z$ , such that the model can be written in the common linear equation form

$$Y \sim \text{Normal}(b_0 + b_1z, \sigma^2)$$

Linear does not mean that the relationship between  $Y$  and  $X$  is linear; it simply means that  $Y$  can be expressed as a linear function of  $X$ .



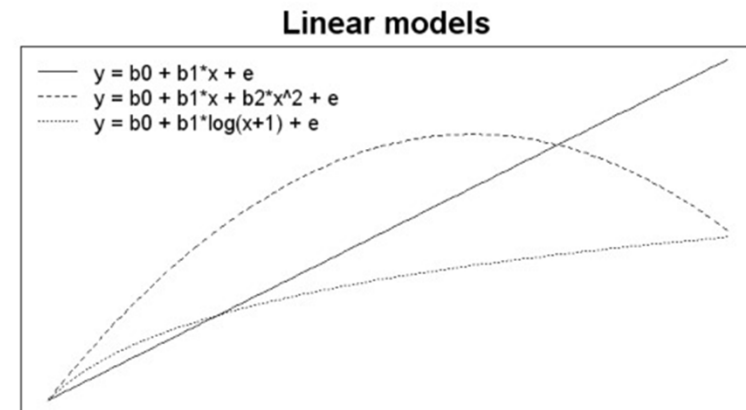
# Linear vs. non-linear models

Example linear models:

$$Y \sim \text{Normal}(b_0 + b_1x, \sigma^2)$$

$$Y \sim \text{Normal}(b_0 + b_1x + b_2x^2, \sigma^2)$$

$$Y \sim \text{Normal}(b_0 + b_1 \log(x), \sigma^2)$$

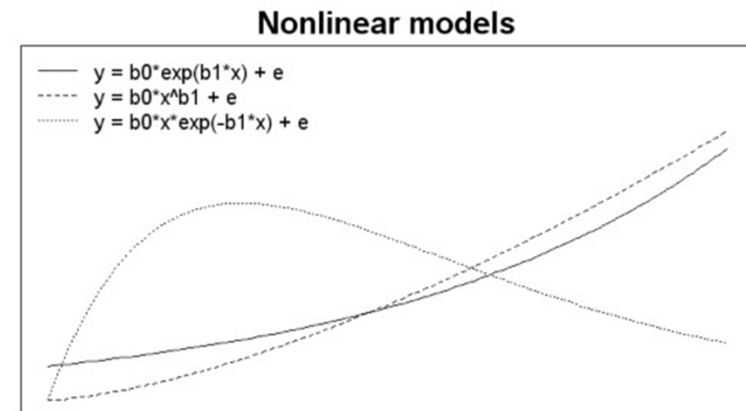


Example nonlinear models:

$$Y \sim \text{Normal}(b_0e^{b_1x}, \sigma^2)$$

$$Y \sim \text{Normal}(b_0x^{b_1}, \sigma^2)$$

$$Y \sim \text{Normal}(b_0xe^{-b_1x}, \sigma^2)$$

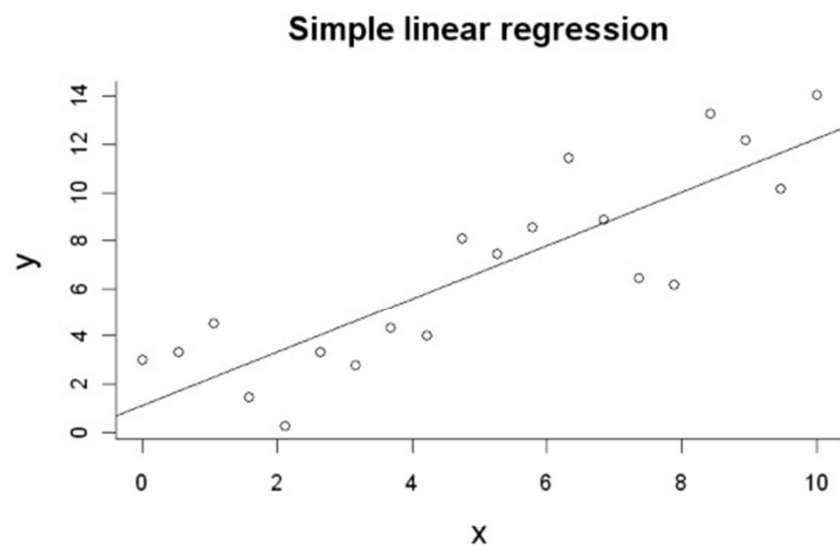


In these examples the models are all linear with respect to  $b_0$  but nonlinear with respect to  $b_1$

## Simple linear regression

- Single continuous predictor

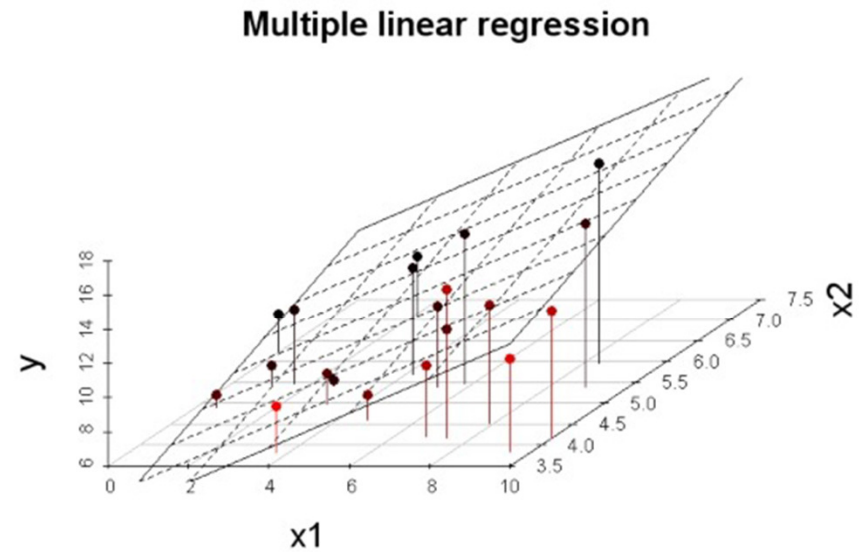
$$Y \sim \text{Normal}(b_0 + b_1x, \sigma^2)$$



## Multiple linear regression

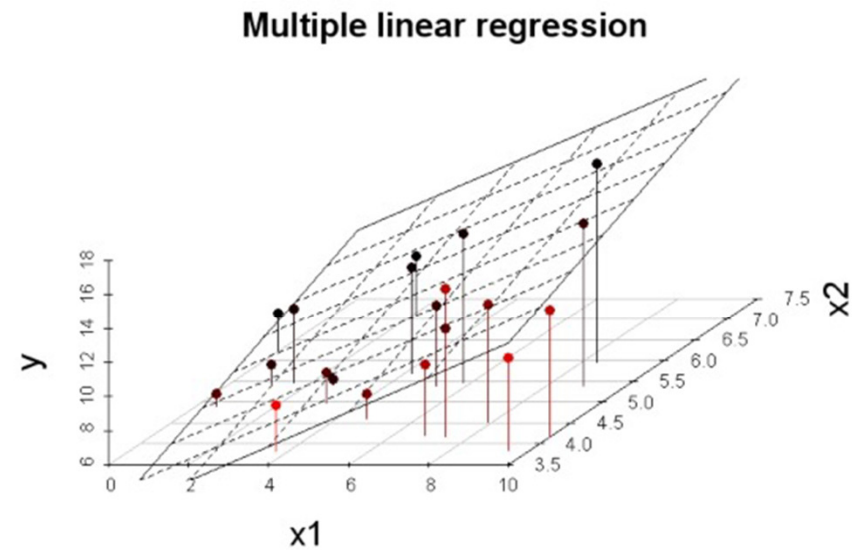
- Multiple continuous predictors

$$Y \sim \text{Normal}(b_0 + b_1x_1 + b_2x_2 + \dots, \sigma^2)$$



## Multiple linear regression

### ■ Multiple continuous predictors



$$Y \sim \text{Normal}(b_0 + b_1x_1 + b_2x_2 + \dots, \sigma^2)$$

In addition, interactions among covariates can be added. This tests whether the slope with respect to one covariate changes linearly as a function of another covariate.

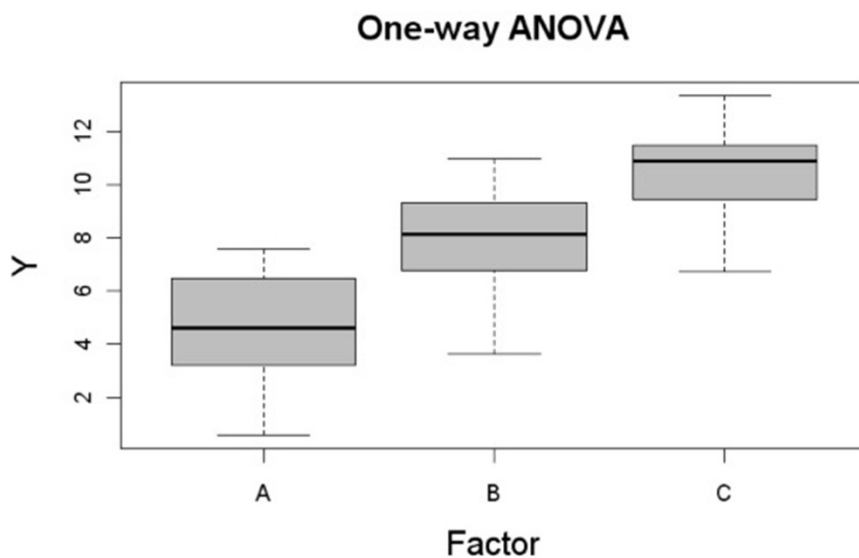
$$Y \sim \text{Normal}(b_0 + b_1x_1 + b_2x_2 + b_{12}x_1x_2, \sigma^2)$$

One-way analysis of variance (ANOVA)

- Single categorical predictor (factor)

$$Y_i \sim \text{Normal}(\alpha_i, \sigma^2)$$

---

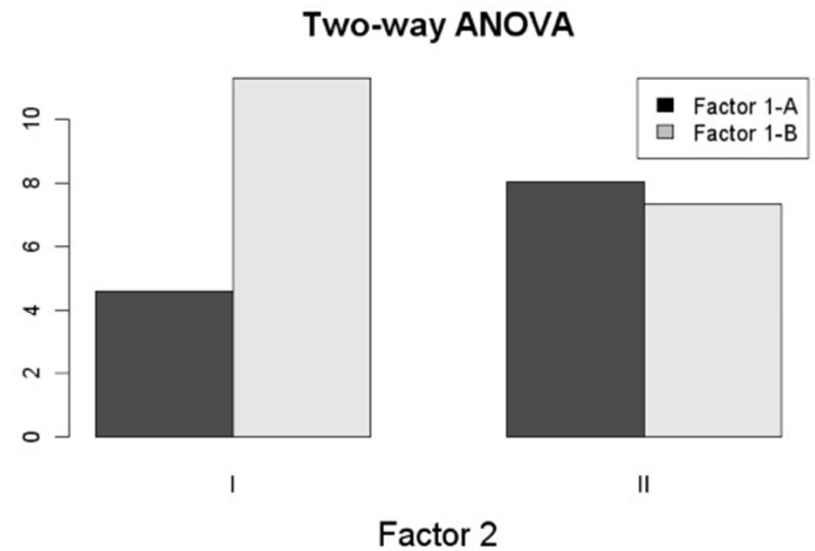


## Multiway ANOVA

- Multiple categorical predictors (factors)

$$Y_{ij} \sim \text{Normal}(\alpha_i + \beta_j + \gamma_{ij}, \sigma^2)$$

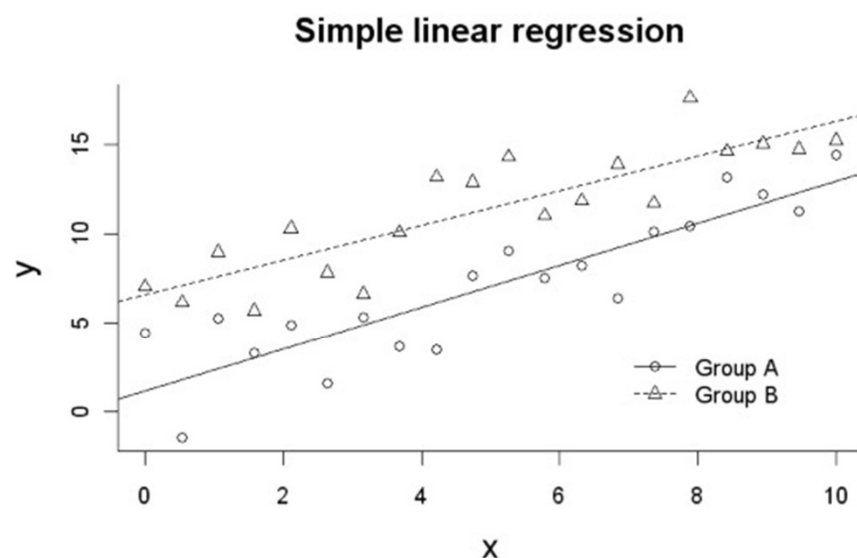
Interaction



## Analysis of covariance (ANCOVA)

- Mix of categorical predictors (factors) and continuous covariate

$$Y_i \sim \text{Normal}(\alpha_i + \beta_i x, \sigma^2)$$



# General linear models

- Linear regression
- One- and multi-way analysis of variance (ANOVA)
- Analysis of covariance (ANCOVA)
- In R these procedures are use the function `lm()`